

News Session-Based Recommendations using Deep Neural Networks

Gabriel de Souza Pereira
Moreira*
CI&T
Campinas, SP, Brazil
gabrielpm@ciandt.com

Felipe Ferreira
Globo.com
Rio de Janeiro, RJ, Brazil
felipe.ferreira@corp.globo.com

Adilson Marques da Cunha
Brazilian Aeronautics Institute of
Technology - ITA
São José dos Campos, SP, Brazil
cunha@ita.br

ABSTRACT

News recommender systems are aimed to personalize users experiences and help them to discover relevant articles from a large and dynamic search space. Therefore, news domain is a challenging scenario for recommendations, due to its sparse user profiling, fast growing number of items, accelerated item's value decay, and users preferences dynamic shift.

Some promising results have been recently achieved by the usage of Deep Learning techniques on Recommender Systems, specially for item's feature extraction and for session-based recommendations with Recurrent Neural Networks.

In this paper, it is proposed an instantiation of the CHAMELEON – a Deep Learning Meta-Architecture for News Recommender Systems. This architecture is composed of two modules, the first responsible to learn news articles representations, based on their text and metadata, and the second module aimed to provide session-based recommendations using Recurrent Neural Networks.

The recommendation task addressed in this work is next-item prediction for users sessions: "what is the next most likely article a user might read in a session?"

Users sessions context is leveraged by the architecture to provide additional information in such extreme cold-start scenario of news recommendation. Users' behavior and item features are both merged in an hybrid recommendation approach.

A temporal offline evaluation method is also proposed as a complementary contribution, for a more realistic evaluation of such task, considering dynamic factors that affect global readership interests like popularity, recency, and seasonality.

Experiments with an extensive number of session-based recommendation methods were performed and the proposed instantiation of CHAMELEON meta-architecture obtained a significant relative improvement in top-n accuracy and ranking metrics (10% on Hit Rate and 13% on MRR) over the best benchmark methods.

*Also with Brazilian Aeronautics Institute of Technology - ITA.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

DLRS 2018, October 6, 2018, Vancouver, BC, Canada

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-6617-5/18/10...\$15.00

<https://doi.org/10.1145/3270323.3270328>

CCS CONCEPTS

• Information systems → Recommender systems; • Computing methodologies → Neural networks;

KEYWORDS

Recommender Systems; Deep Learning; News Recommendation; Session-Based Recommendation; Context-Based Recommendation; Recurrent Neural Networks

ACM Reference Format:

Gabriel de Souza Pereira Moreira, Felipe Ferreira, and Adilson Marques da Cunha. 2018. News Session-Based Recommendations using Deep Neural Networks. In *3rd Workshop on Deep Learning for Recommender Systems (DLRS 2018)*, October 6, 2018, Vancouver, BC, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3270323.3270328>

1 INTRODUCTION

Recommender Systems (RS) have been increasingly popular in assisting users with their choices, thus enhancing their engagement and overall satisfaction with online services [34]. They are an important part of information and e-commerce systems, enabling users to filter through large information and product spaces.

Recommender systems have been researched and applied in online services from different domains, like music [11] [64] [68] (e.g., Spotify, Pandora, Last.fm), videos (e.g. YouTube [16]), people [3] (e.g., Facebook), jobs [5] (e.g., LinkedIn [36], Xing [49]), and research papers [66] [6] (e.g., Docear [7]), among others.

1.1 Deep Learning on Recommender Systems

Deep Learning (DL) [28] [29] [9] [8] is a hot topic in machine learning communities. The uptake of deep learning by RS community was relatively slow, as the topic became popular only in 2016, with the first Deep Learning for Recommender Systems workshop at the ACM RecSys 2016 [26].

Early pioneer work applying used neural networks to RS was done in [58], where a two-layer Restricted Boltzmann Machine (RBM) slightly outperformed Matrix Factorization.

After a winter on RS research using neural networks, Deep Collaborative Filtering was addressed by [67] and [70] using denoising auto-encoders [65]. Deep neural networks have recently been used to learn item features from unstructured data, like text [4], music [64] [68], and images [47] [23].

Recurrent Neural Networks (RNN) possess several properties that make them attractive for sequence modeling of user sessions. In particular, they are capable of incorporating input from past consumption events, allowing to derive a wide range of sequence-to-sequence mappings [18]. After the seminal work of [25], a research

line has emerged on the usage of RNNs on session-based [27] [24] [69] [45] [59] and session-aware [18] [53] [57] recommendations.

1.2 News Recommender Systems

Popular news portals, such as Google News [15], Yahoo! News [62], The New York Times [61], Washington Post [54] [10], among others have gained increasing attention from a massive amount of online news readers.

Online news recommendations have been addressed by researchers in the last years, either using Content-Based Filtering [41] [12] [56] [32] [50], Collaborative Filtering [15] [17], and Hybrid approaches [14] [44] [41] [55] [43] [42] [63] [20].

News domain poses some challenges for Recommender Systems:

- **Sparse user profiling** – the majority of readers are anonymous and they actually read only a few stories from the entire repository. This results in extreme levels of sparsity in the user-item matrix, as users usually have tracked very little information about their past behaviour, if any [41] [43] [17];
- **Fast growing number of items** – hundreds of new stories are added daily in news portals (e.g., over 300 in The New York Times [61]). This intensifies the cold-start problem, as for fresh items you cannot count on lots of interactions before starting to recommend them [17]. For news aggregators, scalability problems may arise, as a high volume of news articles overload the web within limited time span [50];
- **Accelerated decay of item's value** – information value decays over time. This is specially true in the news domain, as most users are interested in fresh information. Thus, each item is expected to have a short shelf life [15]; and
- **Users preferences shift** - news topics of interest are not as stable as in the entertainment domain. Some user interests shift over time, while other long-term interests remain stable [17]. User's current interest in a session may be affected by his context (e.g., location, access time) [17] or by global context (e.g., breaking news or important events) [20].

2 A DEEP LEARNING ARCHITECTURE FOR NEWS SESSION-BASED RECOMMENDATIONS

In [51], it was proposed the CHAMELEON - a Deep Learning Meta-Architecture for News Recommender Systems. A meta-architecture is a reference architecture that collects together decisions relating to an architecture strategy. It might be instantiated as different architectures with similar characteristics that fulfill a common task, in this case, news recommendations.

As shown in Figure 1, CHAMELEON is composed of two complementary modules, with independent life cycles for training and inference: the Article Content Representation (ACR) and the Next-Article Recommendation (NAR) modules.

In this work, the CHAMELEON meta-architecture was instantiated as a concrete architecture, presented in Figure 2. This instantiation of the ACR module used Convolutional Neural Network (CNN) to learn textual features from news articles. In the NAR module, the sequence of clicks from users sessions was modeled by Long Short-Term Memory (LSTM).

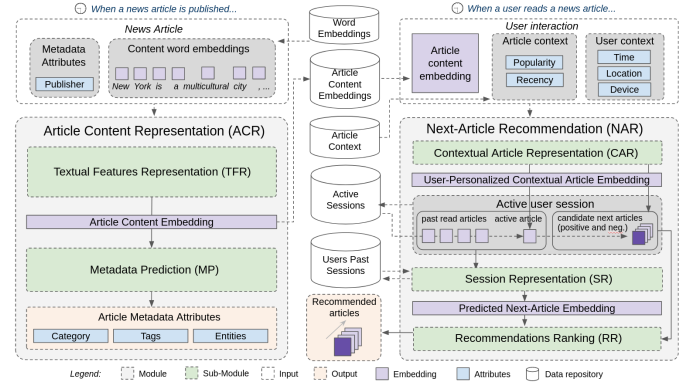


Figure 1: CHAMELEON - A Deep Learning Architecture for News Session-Based Recommendations

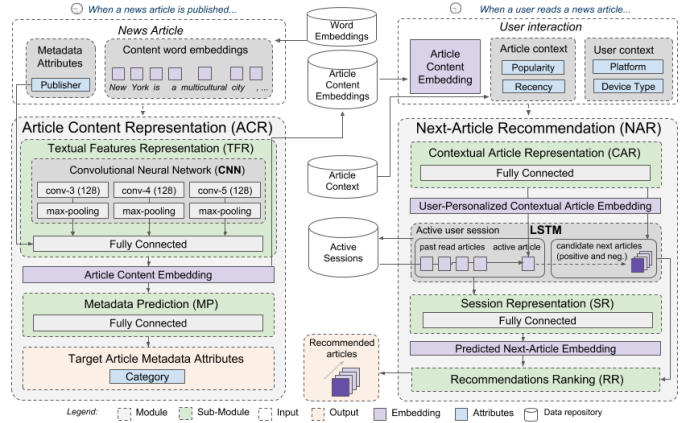


Figure 2: An architecture instantiation of the CHAMELEON, using CNN and LSTM

The following sections describe this architecture instantiation of the CHAMELEON and also experiments performed on a large news dataset, compared to other session-based recommendation methods.

2.1 Article Content Representation (ACR)

The ACR module is responsible to extract features from news articles text and metadata and to learn a distributed representations (embeddings) for each news article context.

The inputs for the ACR module are (1) article metadata attributes (e.g., publisher) and (2) article textual content, represented as a sequence of word embeddings.

A common practice in Deep Natural Language Processing (NLP) is pre-training word embeddings using methods like Word2Vec [48] and GloVe [52] in a larger text corpus of the target language (e.g., Wikipedia).

In this instantiation of the *Textual Features Representation (TFR)* sub-module from ACR module, 1D CNNs were used to extract features from textual items, like in [40] and [13].

Article's textual features and metadata inputs were combined by using a sequence of Fully Connected (FC) layers to produce *Article Content Embeddings*.

For scalability reasons, *Article Content Embeddings* are not directly trained for recommendation task, but for a side task of news metadata classification. In this work, they were trained to classify the category (editorial section) of news articles.

In the last neural network layer, the *softmax* function was used to normalize the output layer as a probability distribution, as follows:

$$\sigma(x_j) = \frac{e^{x_j}}{\sum_i e^{x_i}}. \quad (1)$$

and cross-entropy log loss is used for optimization, as follows:

$$l(\theta) = -\frac{1}{N} \left(\sum_{i=1}^N y_i \cdot \log(\hat{y}_i) \right) + \lambda \|\theta\|_2, \quad (2)$$

where y is a vector with the one-hot encoded label for each instance, \hat{y} is the vector with the output probabilities for each class, previously normalized by *softmax* function, θ , representing model parameters to be learned, and λ to control the importance of the regularization term, to avoid overfitting.

The trained *Article Content Embeddings* are stored in a repository, for further usage by *NAR* module.

2.2 Next-Article Recommendation (NAR)

The *Next-Article Recommendation (NAR)* module is responsible for providing news articles recommendations for active sessions.

Due to the high sparsity of users and their constant interests shift, this work leverages only session-based contextual information, ignoring possible users' past sessions.

The inputs for the *NAR* module are: (1) the pre-trained *Article Content Embedding* of the last viewed article; (2) the contextual properties of the article (popularity and recency); and (3) the user context (e.g. time, location, and device). These inputs are combined by Fully Connected layers to produce a *User-Personalized Contextual Article Embedding*, whose representations might differ for the same article, depending on the user context and on the current article context (popularity and recency).

The *NAR* module uses a type of RNN – the Long Short-Term Memory (LSTM) [30] – to model the sequence of articles read by users in their sessions, represented by their *User-Personalized Contextual Article Embeddings*. For each article of the sequence, the RNN outputs a *Predicted Next-Article Embedding* – the expected representation of a news content the user would like to read next in the active session.

In most deep learning architectures proposed for RS, the neural network outputs a vector whose dimension is the number of available items. Such approach may work for domains where the items number is more stable, like movies and books. Although, in the dynamic scenario of news recommendations, where thousands of news stories are added and removed daily, such approach could require full retrain of the network, as often as new articles are published.

For this reason, instead of using a softmax cross-entropy loss, the *NAR* module is trained to maximize the similarity between the

Predicted Next-Article Embedding and the *User-Personalized Contextual Article Embedding* corresponding to the next article actually read by the user in his session (positive sample), whilst minimizing its similarity with negative samples (articles not read by the user during the session).

With this strategy, a newly published article might be immediately recommended, as soon as its *Article Content Embedding* is trained and added to a repository. The inspiration for this approach came from the DSSM [31] and its derived works for RS, like the MV-DNN [19], the TDSSM [60], and the RA-DSSM [38], which uses a ranking loss based on embeddings similarity.

The *Predicted Next-Article Embedding* and the *User-Personalized Contextual Article Embedding*, further referred as p and $item$, are vectors with the same arbitrary dimension. In Equation 3, it is defined a relevance function R as the cosine similarity between $item$ and p , as shown in Equation 4.

$$R(s, item) = \cos(s, item) \quad (3)$$

$$\cos(\theta) = \frac{a \cdot b}{\|a\| \|b\|} \quad (4)$$

Ranking-based loss functions are usually suitable for Top-N recommendations. The objective of the *NAR* module is to produce a ranked list of the next likely article ($item \in D$, where D is the set of all items) the user will read in the session (next-click prediction). Thus, the model should learn to maximize the similarity between the *Predicted Next-Article Embedding* (p) and the *User-Personalized Contextual Article Embedding* of the next article read by the user ($item^+$), whilst minimizing the pairwise similarity between p and *User-Personalized Contextual Article Embeddings* of negative samples ($item^- \in D^-$), where D^- is the set of all items not read by the user in his session.

As D set may be a very large in news domain, it is approximated as a set D' – the union of the unit set with the clicked item (positive sample) $\{item^+\}$ and a set with random negative samples from D^- .

The posterior probability of the next-clicked article given an active user session was computed by using a *softmax* function over the relevance score proposed by the DSSM [31], as shown in Equation 5,

$$P(item^+ | s) = \frac{\exp(\gamma R(p, item^+))}{\sum_{item \in D'} \exp(\gamma R(p, item))} \quad (5)$$

where γ is a smoothing factor (usually referred to as *temperature*) for the softmax function, which may be a trainable parameter or empirically set on a held-out data set.

The *NAR* module neural network parameters are estimated to maximize the likelihood of the next-clicked article given the user session. The loss function to be minimized, also introduced by the DSSM [31], is shown in Equation 6,

$$l(\theta) = -\log \prod_{(p, item^+)} P(item^+ | s), \quad (6)$$

where θ represent the model parameters to be learned. Since $l(\theta)$ is differentiable w.r.t. to θ , the *NAR* module is trained using back-propagation on gradient-based numerical optimization algorithms.

3 EXPERIMENTS

The proposed instantiation of CHAMELEON meta-architecture was implemented using TensorFlow [1], a popular Deep Learning framework. The source code for this neural architecture and also the baseline methods implemented for these experiments was open-sourced¹.

The neural network models were trained and evaluated on Google Cloud Platform ML Engine, a managed platform for Deep Learning.

For these experiments, a proprietary dataset was provided by the Globo.com². The Globo.com is the most popular news portal in Brazil, with more than 80M unique users and 100k new contents per month. The dataset sample contained user interactions from Oct. 1 to 16, 2017, including more than 3 million clicks, distributed in 1.2 million sessions from 330,000 users who read more than 50,000 different news articles during that period.

3.1 Training and evaluation of the ACR module

The ACR module was used to learn the *Article Content Embeddings* for news articles. The model was trained to classify articles categories (editorial subsection in the news portal) based on its textual content and metadata. The dataset contained 364k news articles from 461 categories. The distribution of the top 200 categories can be seen in Figure 3.

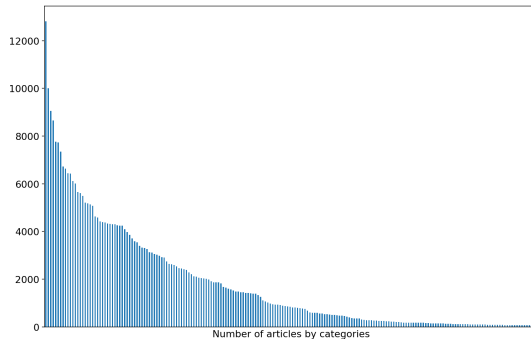


Figure 3: Distribution of number of articles by categories

Articles texts were represented by sequences of pre-trained word embeddings for Portuguese language³. Textual representation was learned by three 1D CNN layers, with window sizes of 3, 4, and 5 (to model word n-grams), each with 128 filters. After max-pooling operation, each of the three layers outputs feature maps (128-dims vectors), which were concatenated with other article's metadata (publisher) by means of a Fully Connected layer.

Training and evaluation were performed using the same dataset, as the objective for this network was not generalization, but to learn representations (*Article Content Embeddings*) for articles content and metadata, with dimension size of 250.

Figure 4 presents a visualization produced using t-SNE, with sampled *Article Content Embeddings* for the 15 categories with most articles. It can be observed that articles are clustered around their

categories, which is expected, since the embeddings were trained to classify articles categories.

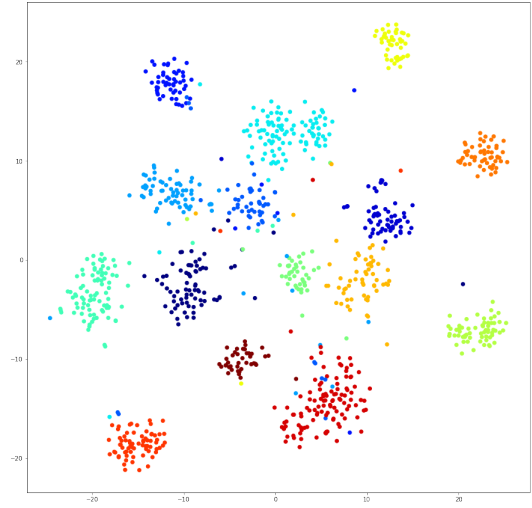


Figure 4: t-SNE visualization of sampled *Article Content Embeddings*, for the top 15 categories

After the training, *Article Content Embeddings* were persisted in a repository, for further usage by NAR module.

3.2 Training and evaluation of the NAR module

In the Globo.com dataset, a session represents a sequence of user clicks with no more than 30 minutes between interactions.

To train NAR module, user interactions sequences were grouped by session and ordered by event time. Sessions with only 1 interaction (invalid for next-click prediction) and with more than 20 interactions (outlier users, or possibly bots) were discarded.

Each training or evaluation mini-batch was composed by a number of sessions (sequences of clicked items), padded to the largest session within the mini-batch.

A cross-validation approach is commonly used to evaluate machine learning models, where samples are split randomly to train or validation sets. Therefore, ignoring the temporal factor of news interactions would be unrealistic.

News readership is very dynamic, as global interests may suddenly shift due to breaking events (e.g., natural disasters, or a royal family member birth) or may follow some seasonality (e.g., soccer during the World Cup, politics during a presidential election, etc.) [14] [22]. The popularity of a news article often changes very quickly, decaying in function of hours.

For this reason, it was devised for NAR module a temporal offline evaluation method, which emulates a real-world scenario of continuously training the model with streaming user clicks and deploying a new trained model once an hour, to provide recommendations for new user sessions. The temporal offline evaluation method is described as follows:

- (1) Train the NAR module with sessions within the active hour; and

¹https://github.com/gabrielspmoreira/chameleon_recsys

²<http://g1.globo.com/>

³It was used a pre-trained Word2Vec *skip-gram* model (300-dims) for Portuguese language, available in <http://nilc.icmc.usp.br/embeddings>

- (2) Evaluate the *NAR* module with sessions within the next hour, for the task of the next-click prediction.

This method is also scalable because, as each session is used only once for *NAR* module training (online learning), there is no need to train on past sessions again (full retrain).

During the training and evaluation loop, in order to keep the consistency of temporal contextual information of articles (popularity and recency), the following method was created:

- (1) Keep a global buffer with the last N clicks (article reads), considering all users;
- (2) Compute articles recent popularity by counting their clicks within the buffer; and
- (3) Compute articles recency as the number of elapsed hours since article was published; and
- (4) For each article read by a user, look up for the updated article context features (recent popularity and recency).

For each clicked article, the features for the *NAR* module was its pre-trained *Article Content Embedding*; the article context attributes (popularity and recency), smoothed by a *log* function; and the user context attributes: platform (web, app) and device type (desktop, mobile, tablet). These inputs were concatenated and normalized with Layer Normalization technique [2].

The *NAR* module requires a number of negative samples to perform training and evaluation. The strategy adopted was to consider as negative samples any article not read within the session, which was read in other sessions from the training/evaluation mini-batch, as proposed in [25]. When there are not enough negative samples within the batch, items are uniformly sampled from a global buffer with the last N clicks. Such approach may help the network to learn to distinguish between the next clicked item (positive item) and other strong negative candidates (articles recently clicked by other users).

For training, 7 negative samples were used for each session, and for evaluation, 50 negative samples. In this setting, each recommender algorithms is expected to rank the set composed by the clicked item and the negative samples.

The Top- N evaluation metrics used in this study were Hit Rate ($HR@5$) [46] which checks whether the clicked item is present in the top-5 ranked items, and Mean Reciprocal Rank ($MRR@5$) [25] [35] [46], a ranking metric, sensitive to the position of clicked item, which assigns higher score at top ranks.

Hyperparameters were tuned for *NAR* module using in a separate period of the same dataset. The LSTM layer had 255 units and *Predicted Next-Article Embeddings* dimension size was 1024. The best training settings were a mini-batch size of 256 sessions, with a learning rate of $1e-3$, a L_2 regularization factor of $1e-4$ and Adam [37] as a gradient based optimizer to learn model parameters. Xavier initialization method [21] was used to initialize model parameters, as hyperbolic tangent (*tanh*) was the mostly used non-linear function in the *NAR* module.

3.3 Baseline methods

For this experiments, an extensive number of session-based recommendation algorithms was used for comparison.

Despite of the simplicity of some of those methods, they usually have competitive accuracy on session-based recommendations, for

keeping up with the dynamic global reading interests and for being efficiently updated over time [33] [35] [46].

- **GRU4Rec** - Seminal neural architecture using RNNs for session-based recommendations [25]. For this experiment, it was used *GRU4Rec* v2 implementation, with the improvements of [24]⁴;
- **Co-occurrent** - Recommends articles commonly viewed together with the last read article, in other user sessions. This algorithm is a simplified version of the association rules technique, with the maximum rule size of two (pairwise item co-occurrences) [35] [46];
- **Sequential Rules (SR)** - A more sophisticated version of association rules, which considers the sequence of clicked items within the session. A rule is created when an item q appeared after an item p in a session, even when other items were viewed between p and q . The rules are weighted by the distance x (number of steps) between p and q in the session with a linear weighting function $w_{SR} = 1/x$ [46];
- **Item-kNN** - Returns most similar items to the last read article, in terms of the cosine similarity between the vector of their sessions, i.e. it is the number of co-occurrences of two items in sessions divided by the square root of the product of the numbers of sessions in which the individual items are occurred. This was the strongest baseline compared to *GRU4Rec* in [25] [46];
- **Vector Multiplication Session-Based kNN (V-SkNN)** - Compares the entire active session with past sessions and find items to be recommended. The comparison emphasizes items more recently clicked within the session, when computing the similarities with past sessions [33] [35] [46];
- **Recently Popular** - Recommends the most viewed articles from the last N clicks buffer; and
- **Content-Based** - For each article read by the user, recommends similar articles based on the cosine similarity of their *Article Content Embeddings*, from the last N clicks buffer.

3.3.1 Notes and improvements on GRU4Rec. The *GRU4Rec* v2 [24] implementation have added the capability of incremental re-training, which supports training on sessions with new items, without the need to retrain on the full dataset. In this setting, it dynamically adds input and output neuron units corresponding to new items, whose connections are randomly initialized.

Therefore, the *GRU4Rec* does not support recommending items not seen on training. For this reason, during its evaluation, it was necessary to ignore fresh articles published since the last training (once each hour), which corresponded to about 2% of ignored clicked items. This approach might have slightly overestimated *GRU4Rec* accuracy, as recommending unknown items tends to be more challenging.

GRU4Rec v2 [24] features an interesting negative sampling strategy, where it can be balanced between uniform sampling (all items have the same probability to be picked) or popularity-based sampling (item sampling probability is proportional to its support).

⁴The *GRU4Rec* v2 [24] was released on Jun 12, 2017 and is available in <https://github.com/hidasib/GRU4Rec>

For these experiments, *GRU4Rec*'s uniform sampling lead to a better accuracy than popularity-based sampling. Therefore, as training evolved, it was possible to observe a decreasing accuracy, as the number of unique items increased, thus, the chance to randomly pick old articles as negative samples, which are not relevant anymore.

For this reason, it was implemented for *GRU4Rec* the same sampling strategy used in *CHAMELEON* during training, which uniformly samples negative items from the last N clicks buffer. With this improvement, it was possible to obtain a relative increase of 13% on $HR@5$ and 18% on $MRR@5$ for *GRU4Rec*.

The *GRU4Rec* is very sensitive to hyperparameter choices, which were tuned in a separate period of the same dataset. The best architecture found for *GRU4Rec* used a hidden layer with 300 units, BPR-max loss function, no item embeddings, and no dropout. It was trained with Adam optimizer (momentum=0), learning rate of $1e-4$, L_2 regularization of $1e-5$, and 200 negative samples (from the last N clicks buffer).

3.4 Results

Two experiments were performed in this study, involving different time periods and evaluation frequency:

- (1) Continuous training and evaluating each five hours, during 15 days (Oct. 1-15, 2017); and
- (2) Continuous training and evaluating each hour, on the subsequent day (Oct. 16, 2017).

$HR@5$ and $MRR@5$ metrics were averaged and reported for each evaluation hour.

To keep results comparable, during the *NAR* module evaluation, it was logged the sampled negative items for each session. Thus, the same negative samples by session were used for baseline methods evaluation.

3.4.1 Experiment 1. For this experiment, recommendation models were trained and evaluated on user sessions from a period of 15 days. Articles context attributes (recent popularity and recency) were also updated over time, to emulate a live environment.

After training on sessions corresponding to five hours, sessions within the subsequent hour were used for evaluation.

Figure 5 presents the evolution of $HR@5$ over time, for the sampled evaluation hours. The distribution of the average $HR@5$ by sampled hours can be seen in Figure 6.

It can be observed in Figure 5 that the proposed *CHAMELEON* instantiation keeps constantly a higher accuracy than the baseline methods over time. As shown in Figure 6, the median $HR@5$ for *CHAMELEON* was 0.72, while the best benchmark (SR) got a median of 0.65, a relative improvement of about 10%.

Methods based on k-Nearest Neighbors and Association Rules presented a competitive accuracy on news session-based recommendations, even higher than *GRU4Rec*, as also observed in [33] [35] [46].

The lower accuracy was obtained by the Content-based method, which may indicate that textual similarity is not the best predictor for the next-read article.

The results are similar for $MRR@5$. Figures 7 and 8 show that the accuracy and ranking quality of recommendations provided by *CHAMELEON* was constantly better than benchmarks methods.

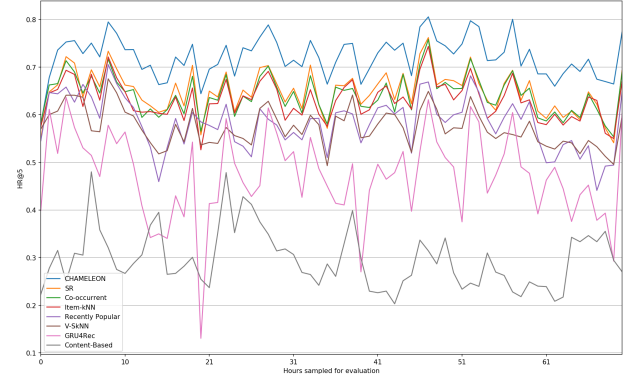


Figure 5: Average $HR@5$ by hour (sampled for evaluation), for a 15-days period

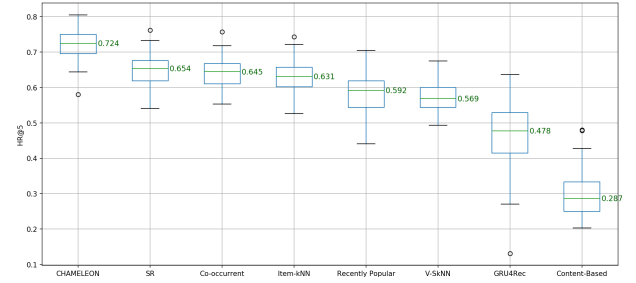


Figure 6: Distribution of average $HR@5$ by hour (sampled for evaluation), for a 15-days period

The median $MRR@5$ for *CHAMELEON* was 0.51, while the best benchmark (SR) got a median of 0.45, a relative improvement of about 13%.

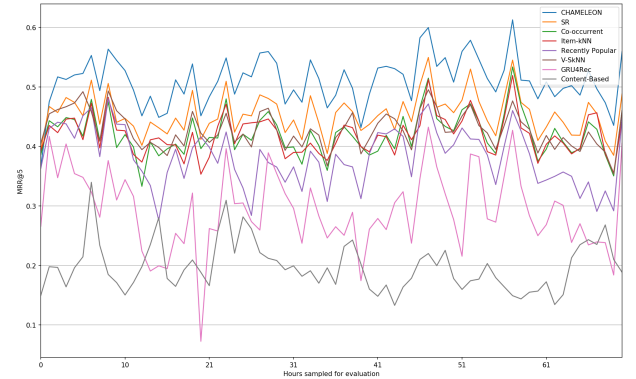


Figure 7: Average $MRR@5$ by hour (sampled for evaluation), for a 15-days period

3.4.2 Experiment 2. In this experiment, recommendation methods were evaluated more often (once an hour), for a period of 24 hours (Oct. 16, 2017). Thus, after training incrementally on sessions from each hour, sessions of the next hour were used for evaluation.

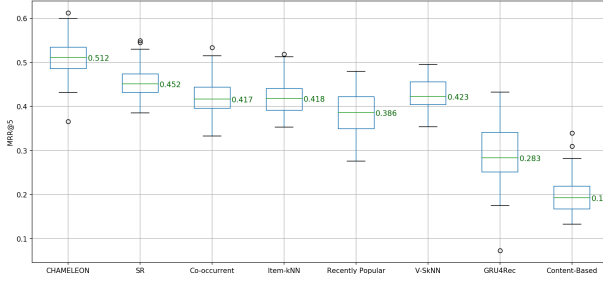


Figure 8: Distribution of average MRR@5 by hour (sampled for evaluation), for a 15-days period

The models trained in *Experiment 1* were used to initialize (parameters and states) models for *Experiment 2*, to emulate a RS that has already being trained for some days.

Figures 9 and 10 presents the evolution of the average HR@5 and MRR@5 by hour, within a period of 24 hours. Once again, it can be seen that the accuracy and ranking quality obtained by the proposed *CHAMELEON* instantiation keeps higher than the baseline methods, throughout the day.

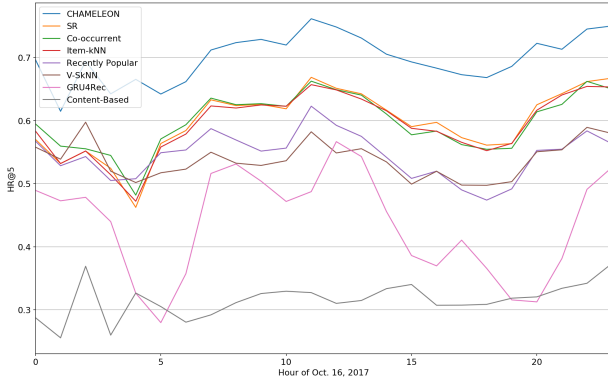


Figure 9: Average HR@5 by hour, for Oct. 16, 2017

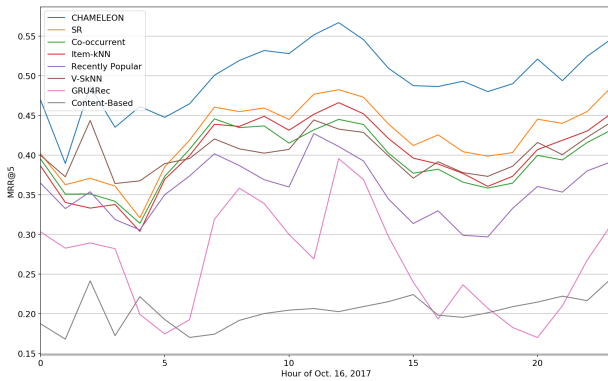


Figure 10: Average MRR@5 by hour, for Oct. 16, 2017

4 RELATED WORK

One of the main inspirations for the *CHAMELEON* was the *GRU4Rec* [25], the seminal work on the usage of Recurrent Neural Networks (RNN) on session-based recommendations, and subsequent work [27] [24].

Other main inspiration came from the *Multi-View Deep Neural Network (MV-DNN)* [19], which adapted *Deep Structured Semantic Model (DSSM)* [31] for the recommendation task.

The *MV-DNN* maps users and items to a latent space, where the cosine similarity between users and their preferred items is maximized. That approach makes it possible to keep the neural network architecture static, rather than adding new units into the output layer for each new item (e.g., published article), as required by *GRU4Rec* (softmax loss function) [25].

The *MV-DNN* was adapted for news recommendation by [60] *Temporal DSSM (TDSSM)* and [38] *Recurrent Attention DSSM (RA-DSSM)*.

Differently from *CHAMELEON*, *TDSSM* [60] did not model user sessions explicitly, and items and users representations are not directly learned from news content and users behaviours.

The *RA-DSSM* [38] represents articles content by using *Doc2Vec* [39] embeddings (unsupervised training), while *CHAMELEON* trains *Article Content Embeddings* to predict news metadata by using supervised learning. The *RA-DSSM* does not use any contextual information about the user and articles, which may limit its accuracy in a extreme cold-start scenario like news RS. Finally, the offline evaluation for that study has used cross-validation based on random sampling to predict the last article read by the user in a session, whilst for this study it was proposed a temporal offline evaluation method, to mimic a more realist scenario.

5 CONCLUSION

In this study, it was proposed an instantiation of the *CHAMELEON* – a Deep Learning Meta-Architecture for News Recommender Systems, using a CNN to extract textual features from news articles and a LSTM layer to model the sequence of clicked items in user sessions.

It was possible to observe that the recommendations accuracy obtained by the proposed *CHAMELEON* instantiation was constantly higher over time than an extensive number of baseline methods for session-based recommendation, including the popular *GRU4Rec*.

A temporal offline evaluation method was also proposed to emulate the dynamics of news readership, where articles context (recent popularity and recency) is constantly changing. Recommender methods are continuously trained on streaming user clicks and the model may be often deployed (e.g. once an hour) to a production environment, in order to serve recommendations for real users.

This is an ongoing research on news recommendation using the *CHAMELEON* Deep Learning Meta-Architecture. Some of next planned steps are: (1) to understand how different input features (e.g., article content, user context, article context) individually impact the recommendations accuracy; (2) to evaluate recommendations of fresh news articles (item cold-start); (3) to explore other approaches to learn better article content embeddings; and (4) to evaluate models from perspectives other than accuracy, like popularity, coverage and serendipity, in different large news datasets.

It is suggested the adaptation of *CHAMELEON* to provide contextual session-based recommendations in other domains, like e-commerce and media.

ACKNOWLEDGMENTS

The authors would like to thank Globo.com for providing context on its challenges for large-scale news recommender systems and for sharing a dataset to make those experiments possible.

REFERENCES

- [1] Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. 2016. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467* (2016).
- [2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).
- [3] Hernan Badenes, Mateo N Bengualid, Jilin Chen, Liang Gou, Eben Haber, Jalal Mahmud, Jeffrey W Nichols, Aditya Pal, Jerald Schoudt, Barton A Smith, et al. 2014. System U: automatically deriving personality traits from social media for people recommendation. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 373–374.
- [4] Trapit Bansal, David Belanger, and Andrew McCallum. 2016. Ask the gru: Multi-task learning for deep text recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 107–114.
- [5] Mathieu Bastian, Matthew Hayes, William Vaughan, Sam Shah, Peter Skomoroch, Hyungjin Kim, Sal Uryasev, and Christopher Lloyd. 2014. LinkedIn skills: large-scale topic extraction and inference. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 1–8.
- [6] Joeran Beel, Stefan Langer, Marcel Genzmehr, Bela Gipp, Corinna Breiteringer, and Andreas Nürnberger. 2013. Research paper recommender system evaluation: A quantitative literature survey. In *Proceedings of the International Workshop on Reproducibility and Replication in Recommender Systems Evaluation*. ACM, 15–22.
- [7] Joeran Beel, Stefan Langer, Marcel Genzmehr, and Andreas Nürnberger. 2013. Introducing Docear’s research paper recommender system. In *Proceedings of the 13th ACM/IEEE-CS joint conference on Digital libraries*. ACM, 459–460.
- [8] Yoshua Bengio et al. 2009. Learning deep architectures for AI. *Foundations and trends® in Machine Learning* 2, 1 (2009), 1–127.
- [9] Yoshua Bengio, Yann LeCun, et al. 2007. Scaling learning algorithms towards AI. *Large-scale kernel machines* 34, 5 (2007), 1–41.
- [10] R. Bilton. 2016. The Washington Post tests personalized “pop-up” newsletters to promote its big stories. <http://www.niemanlab.org/2016/05/the-washington-post-tests-personalized-pop-up-newsletters-to-promote-its-big-stories/>. (May 2016).
- [11] Jiajun Bu, Shulong Tan, Chun Chen, Can Wang, Hao Wu, Lijun Zhang, and Xiaofei He. 2010. Music recommendation by unified hypergraph: combining social media information and music content. In *Proceedings of the international conference on Multimedia*. ACM, 391–400.
- [12] Michel Capelle, Flavius Frasinca, Marnix Moerland, and Frederik Hogenboom. 2012. Semantics-based news recommendation. In *Proceedings of the 2nd international conference on web intelligence, mining and semantics*. ACM, 27.
- [13] Rose Catherine and William Cohen. 2017. TransNets: Learning to Transform for Recommendation. *arXiv preprint arXiv:1704.02298* (2017).
- [14] Wei Chu and Seung-Taek Park. 2009. Personalized recommendation on dynamic content using predictive bilinear models. In *Proceedings of the 18th international conference on World wide web*. ACM, 691–700.
- [15] Abhinandan S Das, Mayur Datar, Ashutosh Garg, and Shyam Rajaram. 2007. Google news personalization: scalable online collaborative filtering. In *Proceedings of the 16th international conference on World Wide Web*. ACM, 271–280.
- [16] James Davidson, Benjamin Liebald, Junjing Liu, Palash Nandy, Taylor Van Fleet, Ullas Gargi, Sujoy Gupta, Yu He, Mike Lambert, Blake Livingston, et al. 2010. The YouTube video recommendation system. In *Proceedings of the fourth ACM conference on Recommender systems*. ACM, 293–296.
- [17] Jorge Díez Peláez, David Martínez Rego, Amparo Alonso Betanzos, Óscar Lucaces Rodríguez, and Antonio Bahamonde Rionda. 2016. Metrical Representation of Readers and Articles in a Digital Newspaper. In *10th ACM Conference on Recommender Systems (RecSys 2016)*. ACM.
- [18] Tim Donkers, Benedikt Loepp, and Jürgen Ziegler. 2017. Sequential User-based Recurrent Neural Network Recommendations. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. ACM, 152–160.
- [19] Ali Mamdouh Elkahky, Yang Song, and Xiaodong He. 2015. A multi-view deep learning approach for cross domain user modeling in recommendation systems. In *Proceedings of the 24th International Conference on World Wide Web. International World Wide Web Conferences Steering Committee*, 278–288.
- [20] Elena Viorica Epure, Benjamin Kille, Jon Espen Ingvaldsen, Rebecca Deneckere, Camille Salinesi, and Sahin Albayrak. 2017. Recommending Personalized News in Short User Sessions. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. ACM, 121–129.
- [21] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*. 249–256.
- [22] Jon Atle Gulla, Cristina Marco, Arne Dag Fidjestøl, Jon Espen Ingvaldsen, and Özlem Özgöbek. 2016. The Intricacies of Time in News Recommendation.. In *UMAP (Extended Proceedings)*.
- [23] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [24] Balázs Hidasi and Alexandros Karatzoglou. 2017. Recurrent Neural Networks with Top-k Gains for Session-based Recommendations. *arXiv preprint arXiv:1706.03847* (2017).
- [25] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2016. Session-based recommendations with recurrent neural networks. In *Proceedings of Forth International Conference on Learning Representations*.
- [26] Balázs Hidasi, Alexandros Karatzoglou, Oren Sar-Shalom, Sander Dieleman, Bracha Shapira, and Domonkos Tikk. 2017. DLRS 2017-Second Workshop on Deep Learning for Recommender Systems. In *Proceedings of the 2st Workshop on Deep Learning for Recommender Systems*. ACM, Vol. 29, 34.
- [27] Balázs Hidasi, Massimo Quadrana, Alexandros Karatzoglou, and Domonkos Tikk. 2016. Parallel recurrent neural network architectures for feature-rich session-based recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. ACM, 241–248.
- [28] Geoffrey E Hinton, Simon Osindero, and Yee-Whye Teh. 2006. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
- [29] Geoffrey E Hinton and Ruslan R Salakhutdinov. 2006. Reducing the dimensionality of data with neural networks. *science* 313, 5786 (2006), 504–507.
- [30] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [31] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning deep structured semantic models for web search using clickthrough data. In *Proceedings of the 22nd ACM international conference on Conference on information & knowledge management*. ACM, 2333–2338.
- [32] Ilija Ilić and Sujoy Roy. 2013. Personalized news recommendation based on implicit feedback. In *Proceedings of the 2013 international news recommender systems workshop and challenge*. ACM, 10–15.
- [33] Dietmar Jannach and Malte Ludewig. 2017. When recurrent neural networks meet the neighborhood for session-based recommendation. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. ACM, 306–310.
- [34] Gawesh Jawaheer, Peter Weller, and Patty Kostkova. 2014. Modeling user preferences in recommender systems: A classification framework for explicit and implicit user feedback. *ACM Transactions on Interactive Intelligent Systems (TiiS)* 4, 2 (2014), 8.
- [35] Michael Jugovac, Dietmar Jannach, and Mozghan Karimi. 2018. StreamingRec: A Framework for Benchmarking Stream-based News Recommenders. (2018), 306–310.
- [36] Krishnamurthy Kothapadi, Benjamin Le, and Ganesh Venkataraman. 2017. Personalized Job Recommendation System at LinkedIn: Practical Challenges and Lessons Learned. In *Proceedings of the Eleventh ACM Conference on Recommender Systems*. ACM, 346–347.
- [37] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [38] Vaibhav Kumar, Dhruv Khattar, Shashank Gupta, Manish Gupta, and Vasudeva Varma. 2017. Deep Neural Architecture for News Recommendation. In *Working Notes of the 8th International Conference of the CLEF Initiative, Dublin, Ireland. CEUR Workshop Proceedings*.
- [39] Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *Proceedings of the 31st International Conference on Machine Learning (ICML-14)*. 1188–1196.
- [40] Hanbit Lee, Yeonchan Ahn, Haejun Lee, Seungdo Ha, and Sang-goo Lee. 2016. Quote Recommendation in Dialogue using Deep Neural Network. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 957–960.
- [41] Lei Li, Dingding Wang, Tao Li, Daniel Knox, and Balaji Padmanabhan. 2011. SCENE: a scalable two-stage personalized news recommendation system. In *Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval*. ACM, 125–134.
- [42] Lei Li, Li Zheng, Fan Yang, and Tao Li. 2014. Modeling and broadening temporal user interest in personalized news recommendation. *Expert Systems with Applications* 41, 7 (2014), 3168–3177.
- [43] Chen Lin, Runquan Xie, Xinjun Guan, Lei Li, and Tao Li. 2014. Personalized news recommendation via implicit social experts. *Information Sciences* 254 (2014), 1–18.

- [44] Jiahui Liu, Peter Dolan, and Elin Rønby Pedersen. 2010. Personalized news recommendation based on click behavior. In *Proceedings of the 15th international conference on Intelligent user interfaces*. ACM, 31–40.
- [45] Qiang Liu, Shu Wu, Diyi Wang, Zhaokang Li, and Liang Wang. 2016. Context-aware sequential recommendation. In *Data Mining (ICDM), 2016 IEEE 16th International Conference on*. IEEE, 1053–1058.
- [46] Malte Ludewig and Dietmar Jannach. 2018. Evaluation of Session-based Recommendation Algorithms. *arXiv preprint arXiv:1803.09587* (2018).
- [47] Julian McAuley, Christopher Targett, Qinfeng Shi, and Anton Van Den Hengel. 2015. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 43–52.
- [48] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*. 3111–3119.
- [49] Sonu K Mishra and Manoj Reddy. 2016. A bottom-up approach to job recommendation system. In *Proceedings of the Recommender Systems Challenge*. ACM, 4.
- [50] Itishree Mohallick and Özlem Özgöbek. 2017. Exploring privacy concerns in news recommender systems. In *Proceedings of the International Conference on Web Intelligence*. ACM, 1054–1061.
- [51] Gabriel de Souza Pereira Moreira. 2017. CHAMELEON - A Deep Learning Meta-Architecture for News Recommender Systems. In *RecSys'18 Doctoral Symposium, Proceedings of the Twelfth ACM Conference on Recommender Systems*. ACM, 421–425.
- [52] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*. 1532–1543.
- [53] Massimo Quadrana, Alexandros Karatzoglou, Balázs Hidasi, and Paolo Cremonesi. 2017. Personalizing Session-based Recommendations with Hierarchical Recurrent Neural Networks. In *Proceedings of the 11th ACM Conference on Recommender Systems*.
- [54] Graff, R. 2015. How the Washington Post used data and natural language processing to get people to read more news. <https://knightlab.northwestern.edu/2015/06/03/how-the-washington-posts-clavis-tool-helps-to-make-news-personal/>. (June 2015).
- [55] Junyang Rao, Aixia Jia, Yansong Feng, and Dongyan Zhao. 2013. Personalized news recommendation using ontologies harvested from the web. In *International Conference on Web-Age Information Management*. Springer, 781–787.
- [56] Hongda Ren and Wei Feng. 2013. Concert: A concept-centric web news recommendation system. In *International Conference on Web-Age Information Management*. Springer, 796–798.
- [57] Massimiliano Ruocco, Ole Steinar Lillestøl Skrede, and Helge Langseth. 2017. Inter-Session Modeling for Session-Based Recommendation. In *Proceedings of the 2nd Workshop on Deep Learning for Recommender Systems*. ACM, 24–31.
- [58] Ruslan Salakhutdinov, Andriy Mnih, and Geoffrey Hinton. 2007. Restricted Boltzmann machines for collaborative filtering. In *Proceedings of the 24th international conference on Machine learning*. ACM, 791–798.
- [59] Elena Smirnova and Flavian Vasile. 2017. Contextual Sequence Modeling for Recommendation with Recurrent Neural Networks. In *Proceedings of the 11th ACM Conference on Recommender Systems*.
- [60] Yang Song, Ali Mamdouh Elkahky, and Xiaodong He. 2016. Multi-rate deep learning for temporal recommendation. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*. ACM, 909–912.
- [61] A. Spangher. 2015. Building the Next New York Times Recommendation Engine. <https://open.blogs.nytimes.com/2015/08/11/building-the-next-new-york-times-recommendation-engine/>. (Aug 2015).
- [62] Michele Trevisiol, Luca Maria Aiello, Rossano Schifanella, and Alejandro Jaimes. 2014. Cold-start news recommendation with domain-dependent browse graph. In *Proceedings of the ACM Recommender System conference, RecSys*, Vol. 14.
- [63] Michele Trevisiol, Luca Maria Aiello, Rossano Schifanella, and Alejandro Jaimes. 2014. Cold-start news recommendation with domain-dependent browse graph. In *Proceedings of the 8th ACM Conference on Recommender systems*. ACM, 81–88.
- [64] Aaron Van den Oord, Sander Dieleman, and Benjamin Schrauwen. 2013. Deep content-based music recommendation. In *Advances in neural information processing systems*. 2643–2651.
- [65] Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. 2008. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*. ACM, 1096–1103.
- [66] Chong Wang and David M Blei. 2011. Collaborative topic modeling for recommending scientific articles. In *Proceedings of the 17th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 448–456.
- [67] Hao Wang, Naiyan Wang, and Dit-Yan Yeung. 2015. Collaborative deep learning for recommender systems. In *Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. ACM, 1235–1244.
- [68] Xinxi Wang and Ye Wang. 2014. Improving content-based and hybrid music recommendation using deep learning. In *Proceedings of the 22nd ACM international conference on Multimedia*. ACM, 627–636.
- [69] Caihua Wu, Junwei Wang, Juntao Liu, and Wenyu Liu. 2016. Recurrent neural network based recommendation for time heterogeneous feedback. *Knowledge-Based Systems* 109 (2016), 90–103.
- [70] Yao Wu, Christopher DuBois, Alice X Zheng, and Martin Ester. 2016. Collaborative denoising auto-encoders for top-n recommender systems. In *Proceedings of the Ninth ACM International Conference on Web Search and Data Mining*. ACM, 153–162.